

## Density Estimation [Non-parametric]



$$\therefore P(\text{points in } R) = \frac{k}{N} = \frac{4}{13}$$

let  $V$  be volume of  $R$ , then for  $\underline{x}$  in  $R$ :

$$P(\underline{x}) = \frac{\frac{k}{N}}{V} \rightarrow \text{prob in a region}$$

$\rightarrow$  area of a region

For better estimation

1.  $V \rightarrow 0$
2.  $k \rightarrow \infty$
3.  $\frac{k}{n} \rightarrow 0$

## Method 1: Kernel Density Estimation (KDE)

- o Create a window or kernel function  $\Delta(\underline{x})$ 
  - e.g. Parzen Window
- o Center the function at each points  $\Delta(\underline{x} - \underline{x}_i)$

- o Take sum of all functions

$$P_n(\underline{x}) = \frac{1}{n} \sum_{i=1}^n \Delta(\underline{x} - \underline{x}_i)$$

We can define  $\Delta(\underline{x} - \underline{x}_i)$  as  $\frac{1}{V_n} K\left(\frac{\underline{x} - \underline{x}_i}{h}\right)$ , choose  $V_n$  to ensure

$$\int \frac{1}{V_n} K\left(\frac{\underline{x} - \underline{x}_i}{h}\right) d\underline{x} = 1.$$

$$\text{Then } P_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{V_h} K\left(\frac{x-x_i}{h}\right)$$

Usually,  $V_h$  could be  $h$ , then:

$h$  is width  
( $x_i$  is within this width)

$$P_h(x) = \frac{1}{nh} \sum_{i=1}^h K\left(\frac{x-x_i}{h}\right)$$

## Method 2: K-Nearest Neighbors (KNN)

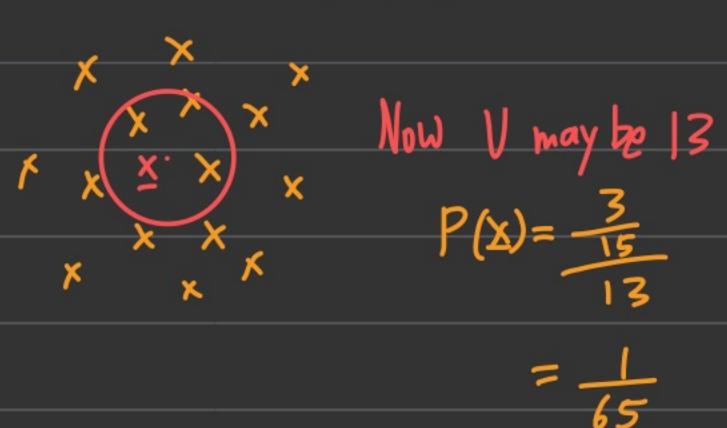
No width ( $h$ ) parameter

No window/kernel

use  $P(x) = \frac{k}{N}$

let  $k=3, N=15$

- Choose  $k$  as a constant
- $N$  is calculated based on  $k$



## Classification with Density Estimation [Generative Approach]

1. Estimate  $P(x|S_k)$  separately for each class  $S_k$ . [only use points in  $S_k$ ]

2. Use or estimate  $P(S_k)$ . e.g.  $\hat{P}(S_k) = \frac{n^{(k)} \rightarrow \text{points\# in } S_k}{N} \rightarrow \text{total points\#}$

3. Use Bayes Classifier. [ $\underbrace{\hat{P}(x|S_k)\hat{P}(S_k)}$  >  $\hat{P}(x|S_j)\hat{P}(S_j)$   $\forall j \neq k$ ]

$$\hat{P}(x|S_k)\hat{P}(S_k) > \hat{P}(x|S_j)\hat{P}(S_j) \quad \forall j \neq k$$

$$\Rightarrow x \in S_k$$

## Parameter Estimation [Classification Case]

Instead of estimate  $P(\underline{x} | S_k)$  with KDE or KNN,

assume  $P(\underline{x} | S_k) = f(\underline{x}, \underline{\theta}_k)$  ↗ known or assumed.  
 ↘ calculated with points in  $S_k$

Ex:  $p(\underline{x} | S_k) = N(\underline{x}, \underline{\mu}_k, \underline{\Sigma}_k)$   
 ↗ unknown

① Consider  $\underline{\theta}$  as deterministic

MLE (Maximum Likelihood Estimation)

Ex:  $p(\underline{x} | \underline{\theta}) = N(\underline{x}, \underline{\mu}, \underline{\Sigma})$  ↗ known  
 ↗ unknown ← multivariate Gaussian

$$\therefore \hat{\underline{\theta}}_{ML} = \arg \max_{\underline{\theta}} p(\underline{x}_1, \dots, \underline{x}_n | \underline{\theta}) \Rightarrow L = \ln p(\underline{x}_1, \dots, \underline{x}_n | \underline{\theta}) = \ln \pi P(\underline{x}_k | \underline{\theta}) = \sum \ln p(\underline{x}_k | \underline{\theta})$$

$$\therefore \hat{\underline{\mu}}_{ML} = \arg \max_{\underline{\mu}} L$$

$$\nabla_{\underline{\mu}} L = \sum_{k=1}^n \ln \frac{1}{\sqrt{(2\pi)^n |\underline{\Sigma}|}} e^{-\frac{1}{2} (\underline{x}_k - \underline{\mu})^\top \underline{\Sigma}^{-1} (\underline{x}_k - \underline{\mu})} = \nabla_{\underline{\mu}} \sum_{k=1}^n \ln \frac{1}{\sqrt{(2\pi)^n |\underline{\Sigma}|}} - \frac{1}{2} (\underline{x}_k^\top \underline{\Sigma}^{-1} \underline{x}_k - \underline{x}_k^\top \underline{\Sigma}^{-1} \underline{\mu} - \underline{\mu}^\top \underline{\Sigma}^{-1} \underline{x}_k + \underline{\mu}^\top \underline{\Sigma}^{-1} \underline{\mu}) = 0$$

$$\Rightarrow \sum_{k=1}^n (-\underline{x}_k \underline{\Sigma}^{-1} - \underline{\Sigma}^{-1} \underline{x}_k + \underline{\Sigma}^{-1} \hat{\underline{\mu}}_{ML} + \hat{\underline{\mu}}_{ML} \underline{\Sigma}^{-1}) = 0$$

$$\sum_{k=1}^n (\underline{\Sigma}^{-1} \underline{x}_k - \underline{\Sigma}^{-1} \hat{\underline{\mu}}_{ML}) = 0$$

$$\boxed{\hat{\underline{\mu}}_{ML} = \frac{1}{n} \sum_{k=1}^n \underline{x}_k}$$

$$\text{Now } P(\underline{x} | S_k) = N(\underline{x}, \hat{\underline{\mu}}_{ML}, \underline{\Sigma})$$

Similarly, if  $\Sigma$  is not known

$$\nabla_{\Sigma} L = \nabla_{\Sigma} \sum_K \left[ \ln \frac{1}{\sqrt{\det(\Sigma)}} - \frac{1}{2} (\underline{x}_K - \underline{\mu})^T \Sigma^{-1} (\underline{x}_K - \underline{\mu}) \right] = 0$$

$$= \sum_K \left[ -\frac{1}{2} \Sigma^{-1} + \frac{1}{2} (\underline{x}_K - \underline{\mu})^T \Sigma^{-1} \Sigma^{-1} (\underline{x}_K - \underline{\mu}) \right] = 0$$

$$\Rightarrow \sum_K \left[ -\hat{\Sigma} + (\underline{x}_K - \underline{\mu})(\underline{x}_K - \underline{\mu})^T \right] = 0$$

$$\sum_{k=1}^n (\underline{x}_K - \underline{\mu})(\underline{x}_K - \underline{\mu})^T = n \cdot \hat{\Sigma}$$

$$\Rightarrow \hat{\Sigma}_{ML} = \frac{1}{n} \sum_{k=1}^n (\underline{x}_K - \underline{\mu})(\underline{x}_K - \underline{\mu})^T \quad \text{where } \underline{\mu} = \hat{\mu}_{ML} \text{ or prior}$$

$$\text{Then } P(\underline{x} | \mathcal{S}_K) = N(\underline{x}, \hat{\mu}_{ML}, \hat{\Sigma}_{ML})$$

$$\text{Note: } \frac{\partial \Sigma^{-1}}{\partial \Sigma} = -\Sigma^{-1} \Sigma^{-1} \quad \textcircled{1}$$

$$\frac{\partial \ln |\Sigma|}{\partial \Sigma} = \Sigma^{-1} \quad \text{if } \Sigma \text{ symmetric} \quad \textcircled{2}$$

$$\frac{\partial A^T \Sigma A}{\partial \Sigma} = A A^T \quad \text{if } \Sigma \text{ symmetric} \quad \textcircled{3}$$

② Consider  $\underline{\theta}$  is random

MAP (Maximum A Posterior estimation)

$P(\underline{\theta} | \underline{z})$ : posterior density of  $\underline{\theta}$ .

$$P(\underline{\theta} | \underline{z}) = \frac{P(\underline{z} | \underline{\theta}) \cdot P(\underline{\theta})}{P(\underline{z})}$$

known or assumed like in MLE

$$L = \ln P(\underline{\theta} | \underline{z}) = \underbrace{\ln P(\underline{z} | \underline{\theta})}_{\text{prior on } \underline{\theta}} + \underbrace{\ln P(\underline{\theta}) - \ln P(\underline{z})}_{\text{constant}}$$

Similarly  $\hat{\underline{\theta}}_{\text{MAP}} = \underset{\underline{\theta}}{\operatorname{argmax}} L$

↳ Ex:  $\underline{\theta}$  is size of tumors

$P(\underline{\theta})$  could be  $N(\underline{\mu}_\theta, \Sigma_\theta)$

according to experts,

$$\underline{\mu}_\theta = 8 \quad \Sigma_\theta = 4$$

With prior,  $\hat{\underline{\theta}}_{\text{MAP}}$  could be solved.

Ex: suppose  $P(\underline{x} | \underline{\theta}) = N(\underline{x}, \underline{\mu}, \Sigma)$

$P(\underline{\mu}) = N(\underline{\mu}, \underline{\mu}_\mu, \Sigma_\mu)$  known, prior

$$L = \underbrace{\ln P(x_1, \dots, x_n | \underline{\theta})}_{\text{same as MLE}} + \underbrace{\ln P(\underline{\theta})}_{\text{prior}} = \left[ \sum_K \ln P(x_K | \underline{\theta}) \right] + \ln P(\underline{\theta}) \leftarrow P(\underline{\mu})$$

$$\begin{aligned} \therefore \nabla_{\underline{\mu}} L &= \nabla_{\underline{\mu}} \left\{ \left[ \sum_K \ln \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} e^{-\frac{1}{2} (\underline{x}_K - \underline{\mu})^T \Sigma^{-1} (\underline{x}_K - \underline{\mu})} \right] + \ln \frac{1}{\sqrt{(2\pi)^D |\Sigma_\mu|}} e^{-\frac{1}{2} (\underline{\mu} - \underline{\mu}_\mu)^T \Sigma_\mu^{-1} (\underline{\mu} - \underline{\mu}_\mu)} \right\} \\ &= \sum_K -\frac{1}{2} (-(\Sigma^{-1})^T \underline{x}_K - \Sigma^{-1} \underline{x}_K + \Sigma^{-1} \underline{\mu} + (\Sigma^{-1})^T \underline{\mu}) - \frac{1}{2} \nabla_{\underline{\mu}} (\underline{\mu}^T \Sigma_\mu^{-1} \underline{\mu} - \underline{\mu}^T \Sigma_\mu^{-1} \underline{\mu}_\mu - \underline{\mu}_\mu^T \Sigma_\mu^{-1} \underline{\mu} \\ &\quad + \underline{\mu}_\mu^T \Sigma_\mu^{-1} \underline{\mu}_\mu) \\ &\stackrel{\Sigma^{-1} \text{ is symmetric}}{=} \sum_K -\frac{1}{2} (-2 \Sigma^{-1} \underline{x}_K + 2 \Sigma^{-1} \underline{\mu}) - \frac{1}{2} (2 \Sigma_\mu^{-1} \underline{\mu} - 2 \Sigma_\mu^{-1} \underline{\mu}_\mu) = 0 \end{aligned}$$

$$\Rightarrow \Sigma^{-1} \cdot \sum_{k=1}^n \underline{x}_K - n \cdot \Sigma^{-1} \cdot \hat{\underline{\mu}}_{\text{MAP}} = \Sigma_\mu^{-1} \cdot \hat{\underline{\mu}}_{\text{MAP}} - \Sigma_\mu^{-1} \cdot \underline{\mu}_\mu$$

$$\Rightarrow (n \cdot \Sigma^{-1} + \Sigma_\mu^{-1}) \cdot \hat{\underline{\mu}}_{\text{MAP}} = \Sigma^{-1} \cdot \sum_{k=1}^n \underline{x}_K + \Sigma_\mu^{-1} \cdot \underline{\mu}_\mu$$

where  $\Sigma^{-1}$  is needed.  
 $(\Sigma^{-1} \text{ could be } \Sigma_\mu^{-1})$

Similarly

$$\nabla_{\Sigma} L = \nabla_{\Sigma} \left\{ \left[ \sum_k \ln \frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(\mathbf{x}_k - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x}_k - \boldsymbol{\mu})} \right] + P(\Sigma) \right\}$$

↙ could be matrix Gaussian  
↖ out of my knowing area

Two approaches for  $\Theta_k$  ↪ Ex. as in MLE

1)  $P(\mathbf{x} | \mathcal{S}_k) = N(\mathbf{x}, \boldsymbol{\mu}, \Sigma)$   $\boldsymbol{\mu} = \hat{\boldsymbol{\mu}}_{MAP}$  (point estimate)

2)  $P(\mathbf{x} | \mathcal{S}_k) = \int N(\mathbf{x}, \boldsymbol{\mu}, \Sigma) \underbrace{P(\boldsymbol{\mu}) d\boldsymbol{\mu}}$

Bayes Classification

Posterior density of  $\boldsymbol{\mu}$

Then, use Bayes Classifier.

## Parameter Estimation [Regression Case]

estimate  $P(y|x) = f(y, x, \theta)$

Ex:  $P(y|x) = N(y, \hat{f}(x, \underline{\omega}), \sigma_y^2)$

$\hat{f}(x, \underline{\omega})$  ↪ mean value       $\sigma_y^2$  ↪ unknown

where  $\hat{f}(x, \underline{\omega})$  could be  $\begin{cases} \text{linear: } \underline{\omega}^\top \underline{x} \\ \text{nonlinear: } \underline{\omega}^\top \phi(\underline{x}) \end{cases}$

## Maximum Likelihood Estimation (MLE)

$$\begin{aligned} P(\text{data}|\underline{\theta}) &= P(\underline{x}, \underline{y}|\underline{\theta}) = P(\underline{y}|\underline{x}, \underline{\theta}) \cdot P(\underline{x}|\underline{\theta}) \\ &= P(\underline{y}|\underline{x}, \underline{\theta}) \cdot P(\underline{x}) \quad \uparrow \text{indep.} \end{aligned}$$

$$\therefore L = \ln [P(\underline{y}|\underline{x}, \underline{\theta}) P(\underline{x})] = \ln P(\underline{y}|\underline{x}, \underline{\theta}) + \ln P(\underline{x})$$

$$\therefore \hat{\underline{\theta}}_{\text{ML}} = \underset{\underline{\theta}}{\operatorname{argmax}} L$$

$$\begin{aligned} L &= \ln P(\underline{y}|\underline{x}, \underline{\theta}) + \ln P(\underline{x}) = \ln \prod_{i=1}^n P(y_i|x_i, \underline{\theta}) + \ln \prod_{i=1}^n P(x_i) \\ &= \sum_{i=1}^n \ln P(y_i|x_i, \underline{\theta}) + \sum_{i=1}^n \ln P(x_i) \end{aligned}$$

$$\nabla_{\underline{\theta}} L = \nabla_{\underline{\theta}} \left( \sum_{i=1}^n \ln P(y_i|x_i, \underline{\theta}) \right) = \nabla_{\underline{\theta}} \left( \sum_{i=1}^n \ln N(y_i, \hat{f}(x_i, \underline{\theta}), \sigma_y^2) \right)$$

$$= \nabla_{\underline{\theta}} \left( \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi}\sigma_y} e^{-\frac{1}{2} \left( \frac{y_i - \hat{f}(x_i, \underline{\theta})}{\sigma_y} \right)^2} \right) = \nabla_{\underline{\theta}} \left( \sum_{i=1}^n -\frac{1}{2} \left( \frac{y_i - \hat{f}(x_i, \underline{\theta})}{\sigma_y} \right)^2 \right)$$

$$= + \sum_{i=1}^n \left( \frac{y_i - \hat{\underline{\theta}}^\top \underline{x}_i}{\sigma_y} \right) \cdot \underline{x}_i / \sigma_y = 0$$

$$\therefore \sum_{i=1}^n (y_i - \hat{\underline{\theta}}^\top \underline{x}_i) \underline{x}_i = 0 = (\underline{y} - \underline{\hat{\theta}}^\top \underline{x})^\top \underline{x} = \underline{y}^\top \underline{x} - \underline{\hat{\theta}}^\top \underline{x}^\top \underline{x} = 0$$



## Maximum A Posterior (MAP)

$$\underline{\theta}_{\text{MAP}} = \arg \max_{\underline{\theta}} L = \arg \max_{\underline{\theta}} \ln P(\underline{\theta} | \text{data})$$

Similarly:  $P(\underline{\theta} | \text{data}) = P(\underline{\theta} | \underline{x}, \underline{y}) = \frac{P(\underline{x}, \underline{y} | \underline{\theta}) P(\underline{\theta})}{P(\underline{x}, \underline{y})} \leftarrow \text{constant for } \underline{\theta}$

$$\begin{aligned} \therefore L &= \ln P(\underline{x}, \underline{y} | \underline{\theta}) P(\underline{\theta}) = \ln P(\underline{y} | \underline{x}, \underline{\theta}) \cdot P(\underline{x} | \underline{\theta}) \cdot P(\underline{\theta}) \\ &= \ln P(\underline{y} | \underline{x}, \underline{\theta}) \cdot P(\underline{\theta}) \end{aligned} \quad \nwarrow \text{indep. \& constant for } \underline{\theta}$$

Similarly  $P(\underline{y} | \underline{x}, \underline{\theta}) = \prod_{i=1}^n P(y_i | \underline{x}_i, \underline{\theta}) = \prod_{i=1}^n P(y_i | \underline{x}_i, \underline{\theta})$

and  $P(y_i | \underline{x}_i, \underline{\theta}) = N(y_i, \hat{f}(\underline{x}_i, \underline{\theta}), \sigma_y^2) \quad \nwarrow \text{assume given}$

and  $P(\underline{\theta}) = N(\underline{\theta}, \underline{M}_{\theta}, \underline{\Sigma}_{\theta}) \quad \nwarrow \text{known (prior)}$

where  $\theta_i \& \theta_j, i \neq j, \text{ are indep. ; so } \underline{\Sigma}_{\theta} = \sigma_{\theta}^2 \cdot \underline{I}$

$$\begin{aligned} \Rightarrow \nabla_{\underline{\theta}} L &= \nabla_{\underline{\theta}} \left( \ln \prod_{i=1}^n \underbrace{N(y_i, \hat{f}(\underline{x}_i, \underline{\theta}), \sigma_y^2)}_{\frac{1}{\sigma_y} e^{-\frac{1}{2} \frac{(y_i - \hat{f}(\underline{x}_i, \underline{\theta}))^2}{\sigma_y^2}}} + \ln \underbrace{P(\underline{\theta})}_{\frac{1}{(2\pi)^{\frac{D}{2}} |\underline{\Sigma}_{\theta}|^{\frac{1}{2}}} e^{-\frac{1}{2} [\underline{\theta} - \underline{M}_{\theta}]^T \underline{\Sigma}_{\theta}^{-1} [\underline{\theta} - \underline{M}_{\theta}]}} \right) \\ &= + \sum_{i=1}^n \left( \frac{y_i - \hat{\theta}^T \underline{x}_i}{\sigma_y} \right) \cdot \underline{x}_i / \sigma_y + \nabla_{\underline{\theta}} \left( \ln \frac{1}{(2\pi)^{\frac{D}{2}} |\underline{\Sigma}_{\theta}|^{\frac{1}{2}}} e^{-\frac{1}{2} [\underline{\theta} - \underline{M}_{\theta}]^T \underline{\Sigma}_{\theta}^{-1} [\underline{\theta} - \underline{M}_{\theta}]} \right) \\ &= \frac{1}{\sigma_y^2} \sum_{i=1}^n (y_i - \hat{\theta}^T \underline{x}_i) \cdot \underline{x}_i + -\frac{1}{2\sigma_y^2} \left[ \hat{\theta}^T \hat{\theta} - \hat{\theta}^T \underline{M}_{\theta} - \underline{M}_{\theta}^T \hat{\theta} - \underline{M}_{\theta}^T \underline{M}_{\theta} \right] \\ &= 0 \end{aligned}$$

$$\Rightarrow \frac{1}{\sigma_\theta^2} (\underline{\theta} - \underline{m}_\theta) = \frac{1}{\sigma_y^2} \sum_{i=1}^n (y_i - \underline{\theta}^T \underline{x}_i) \cdot \underline{x}_i = \frac{1}{\sigma_y^2} [(\underline{y} - \underline{\underline{x}} \underline{\hat{\theta}})^T \underline{\underline{x}}]^T = \frac{1}{\sigma_y^2} \underline{\underline{x}}^T (\underline{y} - \underline{\underline{x}} \underline{\hat{\theta}})$$

$$\Rightarrow \frac{1}{\sigma_\theta^2} \underline{\hat{\theta}} + \frac{1}{\sigma_y^2} \underline{\underline{x}}^T \underline{\underline{x}} \underline{\hat{\theta}} = \frac{1}{\sigma_\theta^2} \underline{m}_\theta + \frac{1}{\sigma_y^2} \underline{\underline{x}}^T \underline{y}$$

$$\left( \frac{\sigma_y^2}{\sigma_\theta^2} \cdot \underline{\underline{I}} + \underline{\underline{x}}^T \underline{\underline{x}} \right) \underline{\hat{\theta}} = \frac{\sigma_y^2}{\sigma_\theta^2} \underline{m}_\theta + \underline{\underline{x}}^T \underline{y}$$

$$\Rightarrow \underline{\hat{\omega}}_{MAP} = \underline{\hat{\theta}}_{MAP} = (\underline{\underline{x}}^T \underline{\underline{x}} + \frac{\sigma_y^2}{\sigma_w^2} \underline{\underline{I}})^{-1} (\underline{\underline{x}}^T \underline{y} + \frac{\sigma_y^2}{\sigma_w^2} \underline{m}_w)$$

Similarly, Two approaches for  $\underline{w}$

$$1) P(y | \underline{x}) = N(y, f(\underline{x}, \underline{\hat{w}}_{MAP}), \sigma_y^2) \quad (\text{point estimate})$$

$$2) P(y | \underline{x}) = \int N(y, f(\underline{x}, \underline{w}), \sigma_y^2) \underbrace{P(\underline{w}) d\underline{w}}$$

Bayes Linear Regression ↑

↑ need to find by  $P(\underline{w} | \text{data})$

Predict  $\hat{y}(\underline{x})$

$$\hat{y}(\underline{x}) = E[y | \underline{x}] = \int y P(y | \underline{x}) dy$$

$$\text{We have } P(y | \underline{x}) = N(y, \underbrace{f(\underline{x}, \underline{w})}_{\leftarrow \text{expectation}}, \sigma_y^2)$$

$$\therefore \hat{y}(\underline{x}) = f(\underline{x}, \underline{\hat{w}}) = \underline{\hat{w}}^T \underline{x} < \begin{matrix} \underline{\hat{w}}_{ML} \underline{x} \\ \underline{\hat{w}}_{MAP} \underline{x} \end{matrix}$$

# Bayesian Linear Regression

In MAP, we maximize  $p(\underline{\theta} | \text{data})$  and find  $\hat{\theta}_{\text{MAP}}$

Now, we integrate over it  $\leftarrow P(\underline{\theta} | \text{data})$

$$P(y | \underline{x}) = \int N(y, f(\underline{x}, \underline{w}), \sigma_y^2) P(\underline{w}) d\underline{w}$$

$$\begin{aligned} P(\underline{\theta} | \text{data}) &= P(\underline{\theta} | \underline{x}, \underline{y}) = \frac{P(\underline{x}, \underline{y} | \underline{\theta}) P(\underline{\theta})}{P(\underline{x}, \underline{y})} = \frac{P(\underline{y} | \underline{x}, \underline{\theta}) \cdot P(\underline{x} | \underline{\theta}) \cdot P(\underline{\theta})}{P(\underline{x}, \underline{y})} \\ &= K \cdot P(\underline{y} | \underline{x}, \underline{\theta}) \cdot P(\underline{\theta}) \end{aligned}$$

because data is  
 independent of  $\underline{\theta}$   
 } constant for  $\underline{\theta}$

Similarly,  $P(y_i | \underline{x}_i, \underline{\theta}) = N(y_i, \hat{f}(\underline{x}_i, \underline{\theta}), \sigma_y^2)$   $\leftarrow$  assume given

$$P(\underline{\theta}) = N(\underline{\theta}, \underline{m}_\theta, \underline{\Sigma}_\theta)$$

$\curvearrowright$  known (prior)

where  $\theta_i$  &  $\theta_j$ ,  $i \neq j$ , are indep.; so  $\underline{\Sigma}_\theta = \sigma_\theta^2 \cdot \underline{I}$

$$\therefore P(\underline{y} | \underline{x}, \underline{\theta}) = \prod_{i=1}^n P(y_i | \underline{x}_i, \underline{\theta}) = \prod_{i=1}^n P(y_i | \underline{x}_i, \underline{\theta}) = \prod_{i=1}^n N(y_i, \hat{f}(\underline{x}_i, \underline{\theta}), \sigma_y^2)$$

$$= N(\underline{y}, \hat{f}(\underline{x}, \underline{\theta}), \underline{\Sigma}_y) \quad \text{with } \underline{\Sigma}_y = \sigma_y^2 \underline{I}$$

Then  $P(\underline{\theta} | \text{data}) = K \cdot N(\underline{y}, \hat{f}(\underline{x}, \underline{\theta}), \underline{\Sigma}_y) \cdot N(\underline{\theta}, \underline{m}_\theta, \underline{\Sigma}_\theta)$

Omitting the algebra:

$$P(\underline{\theta} | \text{data}) = N(\underline{\theta}, \underline{m}_p, \underline{\Sigma}_p)$$

$$\text{with } \underline{m}_P = \underline{\Sigma}_P^{-1} (\underline{\Sigma}_{\theta}^{-1} \underline{m}_{\theta} + \frac{1}{\sigma_y^2} \underline{x}^T \underline{y}) = \underline{\Sigma}_P^{-1} \left( \frac{1}{\underline{\Sigma}_{\theta}^{-2}} \underline{m}_{\theta} + \frac{1}{\sigma_y^2} \underline{x}^T \underline{y} \right)$$

$$\underline{\Sigma}_P^{-1} = \underline{\Sigma}_{\theta}^{-1} + \frac{1}{\sigma_y^2} \underline{x}^T \underline{x} = \frac{1}{\underline{\Sigma}_{\theta}^{-2}} \underline{\mathbb{I}} + \frac{1}{\sigma_y^2} \underline{x}^T \underline{x}$$

↑  
Sample correlation of  
feature with output.

↑  
Sample correlation of  
feature with feature.

$$\begin{aligned} \therefore P(y|\underline{x}) &= \int N(y, f(\underline{x}, \underline{\theta}), \sigma_y^2) p(\underline{\theta} | \text{data}) d\underline{\theta} \\ &= \int N(y, f(\underline{x}, \underline{\theta}), \sigma_y^2) N(\underline{\theta}, \underline{m}_P, \underline{\Sigma}_P) d\underline{\theta} \end{aligned}$$

After significant algebra and integration

$$\therefore P(y|\underline{x}) = N(y, \underline{m}_P^T \underline{x}, \sigma^2(\underline{x}))$$

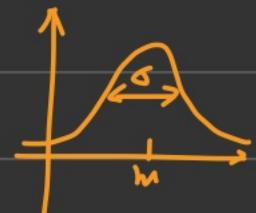
with  $\sigma^2(\underline{x}) = \sigma_y^2 + \underline{x}^T \underline{\Sigma}_P \underline{x}$  and  $\underline{m}_P, \underline{\Sigma}_P$  as before.

$$\therefore \hat{\underline{\omega}} = \underline{m}_P \quad \hat{y}(\underline{x}) = \underline{m}_P^T \cdot \underline{x}$$

## Nonlinear Regression [parameter estimation]

for all formula above:  $\underline{x} \rightarrow \phi(\underline{x}), \underline{\Sigma} \rightarrow \underline{\Phi}$

Ex. We can just set  $\underline{\phi}(\underline{x}) = e^{-\frac{(x-m)^2}{2\sigma^2}}$   $m$ : center  
 $\sigma$ : width



# Summary of parameter estimation

Method	Classification	Regression
Theoretically optimal (know all pdf's)	Bayes min-error (or min-risk)	$\min [TMSE]$ $= \min [E\{\text{square error}\}]$
$\underline{\theta} = \hat{\theta}_{ML}$	✓	✓
$\underline{\theta} = \hat{\theta}_{MAP}$ (point estimate)	✓	✓
Bayesian classification regression (integrate over $\underline{\theta}$ using $P(\underline{\theta}   \text{data})$ )	(not mentioned)	✓